# Multi-Institution Testbed for
# Scalable Digital Archiving

**NSF CISE/Library of Congress DIGARCH Award**



**Stephen Miller**
Scripps Institution of Oceanography

**Bob Detrick**
Woods Hole Oceanographic Institution

**John Helly**
San Diego Supercomputer Center



dg.o 2005 THE NATIONAL CONFERENCE ON DIGITAL GOVERNMENT RESEARCH

**Atlanta 2005-05-17**

# SIOExplorer

## 1. Community Goals

## 2. Barriers to advances
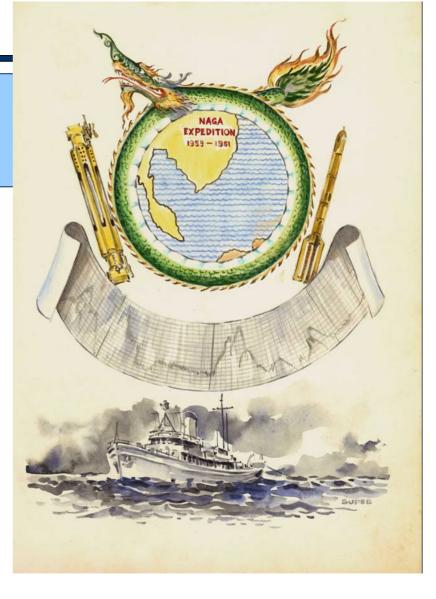
## 3. Cyber-capabilities

SDSC

# 1. Community Goals

**Broad support**

Across disciplines

And institutions

Research

And education

# Guarantee long term preservation



Gulf of California 1939 Expedition, R/V E W Scripps

# Need more than data storage



## Need metadata
Enable re-use

## Also need infrastructure
Networked community tools, archives, understanding

# Why re-use data?

New ship time expensive ($22K/day)

Use archives for:

1. Regional synthesis projects

3. Support other disciplines

3. Monitor environmental changes through time
Before and after
Earthquakes, slumps, seeps
**Volcanoes …**
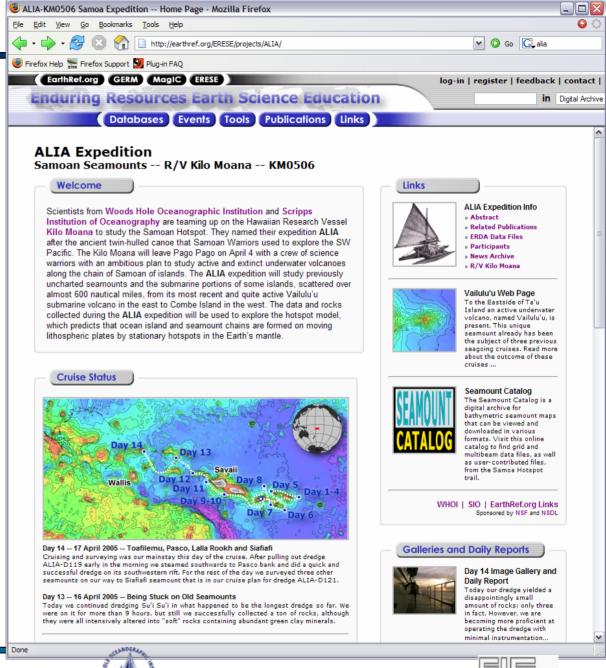
## April 16, 2005

## New Volcanic Cone in the Vailulu'u Crater

With a minimum rate of **eight inches per day**, a new cone has been growing inside the crater of Vailulu'u seamount since the last depth soundings by the US Coastguard vessel Polar Sea in April 2001.

Our survey using the SIMRAD 120 system of the Kilo Moana displays a new volcanic summit at 708 m depth.

This volcano was named Nafanua, after the Samoan Goddess of War.

**Hubert Staudigel, Stan Hart, John Helly, Anthony Koppers, Jasper Kontor**

# 2. Barriers to advances

# Data from a firehose

## Can we keep up?

Shipboard data rates – **yes**

Satellite links – **maybe**
        depends on heading

Metadata – **yes, but**
        not widely implemented



Tiffany Houghton, SDSC, on R/V Sproul

Preservation – **maybe**

Community usage –**help needed from Cyberinfrastructure**

# We can archive from paper documents

Track plots

Cruise reports

Handwritten
and printed data

REVELLE ERA CRUISES
(unedited older data)

JUL 16 1967

| Time | Depth | Remark; | | MAG |
|------|-------|---------|---|-----|
| 0110 | 1570 | S/O 184°TG, 165 RPM | | |
| 15 | 1730 | WS 34 kts, RWD 345° | | 43357 |
| 20 | 1734 | | | |
| 25 | 1579 | | | |
| 30 | 1538 | | | 43360 |
| 35 | 1499 | | | |
| 40 | 1431 | | | |

# But digital preservation is risky business

Endangered Species
9-track tapes →

Exabytes fail

Even CDs fail

RAIDS fail

**"Shoe-box" archiving not to be trusted**

# Solution:  Active Archiving

"Don't trust any media, person or process"

Actively monitor status

Migrate to new storage media

Mirror on multiple systems daily

Backup to independent sites

**Technology makes this possible, just need to do it**

# Example of early backup

Capital burned August 19, 1814

Library of Congress offsite backup

Thomas Jefferson's Library

13th CONGRESS.]    No. 372.    [3d SESSION.

## PURCHASE OF THE LIBRARY OF THOMAS JEFFERSON.

COMMUNICATED TO THE SENATE, OCTOBER 7, 1814.

IN SENATE OF THE UNITED STATES, *October* 7, 1814.

Mr. GOLDSBOROUGH, from the joint committee on the library of Congress, reported:

That they have received, through Mr. Samuel H. Smith, an offer from Mr. Jefferson, late President of the United States, of the whole of his library for Congress, in such a mode, and upon such terms, as they consider highly advantageous to the nation, and worthy the distinguished gentleman who tenders it. But the means placed at the disposal of the committee being very limited and totally inadequate to the purchase of such a library as that now offered, the committee must have recourse to Congress, either to extend their powers, or to adopt such other plan as they may think most proper.

Should it be the sense of Congress to confide this matter to the committee, they respectfully submit the following resolution:

*Resolved, by the Senate and House of Representatives of the United States of America in Congress Assembled,* That the joint Library Committee of the two Houses of Congress be, and they are hereby, authorized and empowered to contract, on their part, for the purchase of the library of Mr. Jefferson, late President of the United States, for the use of both Houses of Congress.

OCTOBER 3, 1814.

SIR:

I have the honor, in furtherance of the proposition contained in a letter from Mr. Jefferson to me, tendering the disposition of his library to Congress, to enclose his letters for submission to the honorable committee over which you preside, with the expression of my readiness at any time to proceed in the discharge of the agency confided to me.

I am, very respectfully, your obedient servant,

SAMUEL H. SMITH.

Hon. ROBERT H. GOLDSBOROUGH,
*Chairman of the Library Committee of Congress.*

MONTICELLO, *September* 21, 1814.

DEAR SIR:

I learn from the newspapers that the vandalism of our enemy has triumphed at Washington over science as well as the arts, by the destruction of the public library, with the noble edifice in which it was deposited.

# 3. Emerging Cyber-capabilities

## SIOExplorer digital library
Design for scalability
Automate harvesting
Collection Builder's Toolkit for other projects

## Crossing institutional boundaries
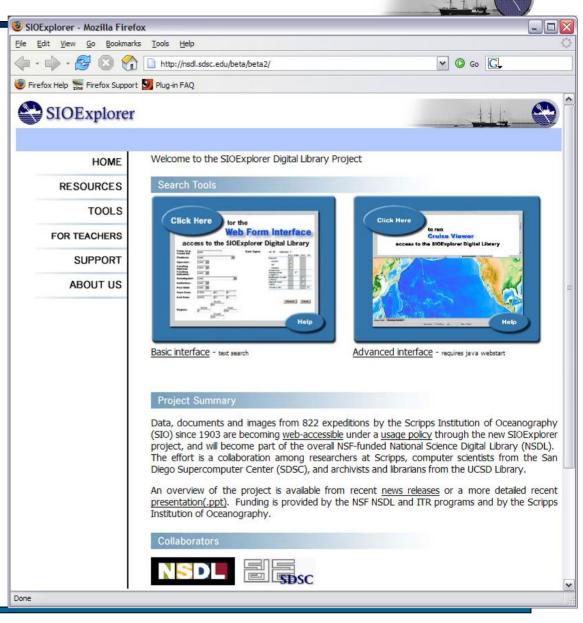Multi-Institution Testbed
SIO, WHOI, SDSC

# SIOExplorer Digital Library

Community access
- Data
- Images
- Documents

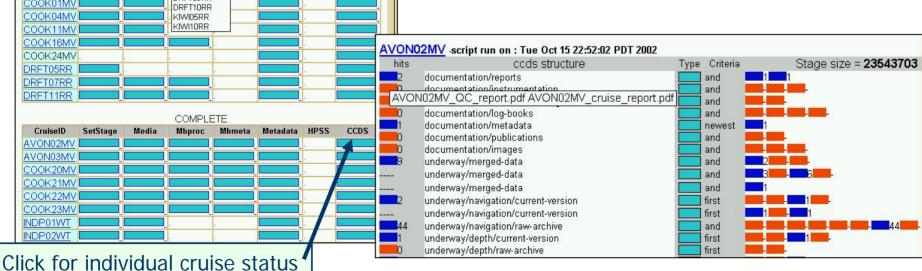647 cruises
150,000 objects
500 GB

Multiple federated collections



**http://SIOExplorer.ucsd.edu**

# Collection status board

Live on web
Auto-updated



Monitor status of 800 cruises, work in progress
4000 files, 10 GB per cruise



Click for individual cruise status

# Issue for future use:
# Access to complete cruise collections

Current practice hit-or-miss

Only selected data streams archived

Cyberinfrastructure allows comprehensive solution

Auto-harvesting and archiving
Data and metadata

**Claim:**
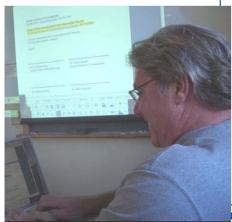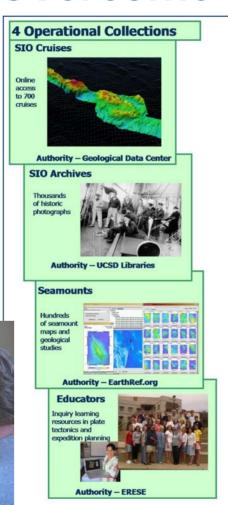**Very little additional cost to archive everything**
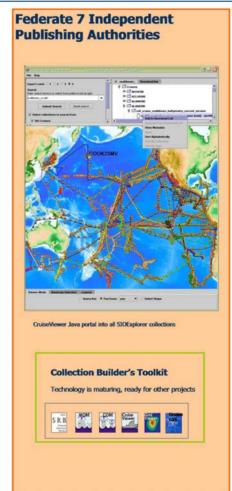
# Design to Overcome Project Barriers

Build scalable digital library

Federate independent authorities

4 Operational collections
3 Work-in-progress



**4 Operational Collections**

**SIO Cruises**

Online access to 700 cruises

*Authority – Geological Data Center*

**SIO Archives**

Thousands of historic photographs

*Authority – UCSD Libraries*

**Seamounts**

Hundreds of seamount maps and geological studies

*Authority – EarthRef.org*

**Educators**

Inquiry learning resources in plate tectonics and expedition planning

*Authority – ERESE*

**Federate 7 Independent Publishing Authorities**

CruiseViewer Java portal into all SIOExplorer collections

**Collection Builder's Toolkit**

Technology is maturing, ready for other projects

S.R.B.    MORE    COBE    Cruise Viewer    Cruise Viewer    Cruise GIS
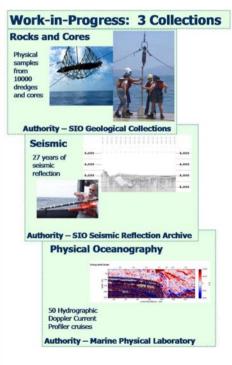
**Investigators and Initial Award**
PI: Brian E. C. Schottlaender, University Librarian, UCSD
Co-PIs: Stephen Miller, Catherine Johnson, Hubert Staudigel, Scripps Institution of Oceanography; John Helly, San Diego Supercomputer Center
NSDL Collections Track 01-21684, Bridging the Gap between Libraries and Data Archives, with additional support from ITR, OCE, ATM and UCSD funds.
Website http://SIOExplorer.ucsd.edu

**Work-in-Progress: 3 Collections**

**Rocks and Cores**

Physical samples from 10000 dredges and cores

*Authority – SIO Geological Collections*

**Seismic**

27 years of seismic reflection

*Authority – SIO Seismic Reflection Archive*

**Physical Oceanography**

50 Hydrographic Doppler Current Profiler cruises

*Authority – Marine Physical Laboratory*

# Multiple access methods

## Google
No interface
Just type name of cruise

## Basic web form
Text-based search for experts

## Java CruiseViewer
Full graphical search

## Web services
Computer-to-computer
Enable next generation interoperability

SDSC

# Java CruiseViewer

Full graphical search
>    All capabilities
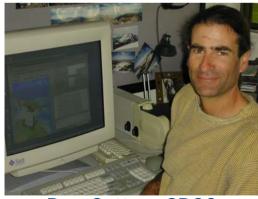>    Any combination of collections
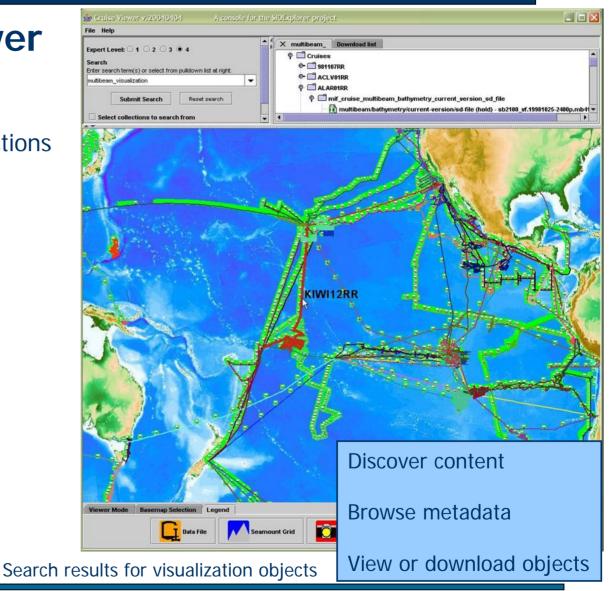
Metadata
>    Oracle or PostgreSQL

Data
>    Storage Resource Broker

User
>    Graphical search
>    Keyword search

Don Sutton, SDSC

Search results for visualization objects

Discover content
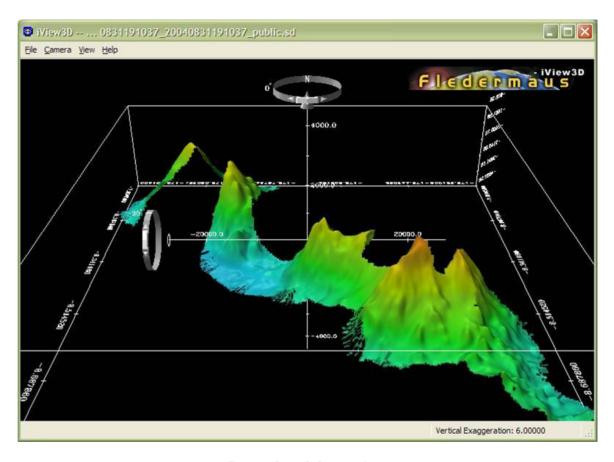
Browse metadata

View or download objects

# Launch visualization experiences

Visualization of multibeam seafloor mapping swath sonar data

300 cruises since 1982
20-km wide swaths

Sonar quality control
Geological research
Education



Download free viewer

http://www.ivs.unb.ca/products/iview3d/

# Broader Impact with ERESE National Teachers Workshops

**Enduring Resources for Earth Science Education**

Two-week summer workshops
        2004 and 2005

Build inquiry-driven learning experiences

**NSDL**
THE NATIONAL SCIENCE DIGITAL LIBRARY

SDSC

# Other organizations using mtf technology

**CUAHSI** Consortium of Universities for Advanced Hydrologic Science, Inc.
Major technology co-development
95 institutional members

**WHOI** – DIGARCH Multi-Institution Testbed project
Bob Detrick

**CCOM/UNH** cruise and multibeam archives
Jim Case, Larry Mayer

**MBARI** – Monterey Bay Aquarium Research Institute collection building in progress
Dave Caress, Andrew Chase

**SOEST/HAWAII** – April 4-26, 2005 realtime digital library testing R/V Kilo Moana

**NIWA** – Digital-Library-in-a-Box tested on R/V Tangaroa in New Zealand
John Helly, Don Robertson

**Arctic DMS** - Data Management System under development
Margo Edwards (Hawaii), Dawn Wright (Oregon State)

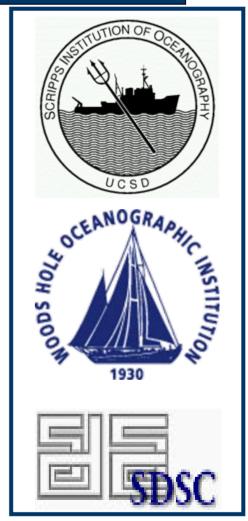# Multi-Institution Testbed for Scalable Digital Archiving

Extend SIOExplorer approach to WHOI

Integrate SIO, SDSC and WHOI tools and data

30 years of WHOI cruise data

4098 Alvin submersible dives

Jason ROV surveys (200 DVD per cruise)

**Results from 1600 NSF awards online**

# Project Challenges

**Auto-harvest data, metadata**
"Shoe-box archives" only
prior to 2002

**Build distributed digital library**
Both institutions
Ships and submersibles

**Extend WHOI data exploration tools**
Persistent digital library objects
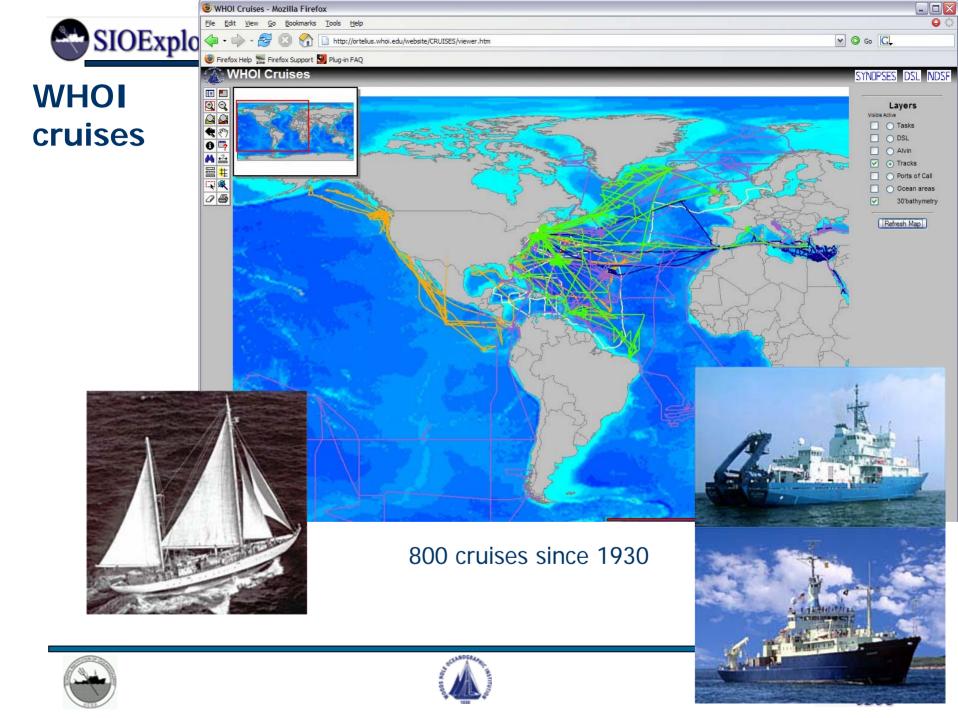
**Interoperability across institutions**



Alvin Frame-Grabber User Interface

# WHOI cruises



800 cruises since 1930

SIOExplo

**4098 Alvin dives**

Since June 26, 1964